

# PIVOTAL ESTIMATION IN HIGH-DIMENSIONAL REGRESSION VIA LINEAR PROGRAMMING

ERIC GAUTIER AND ALEXANDRE B. TSYBAKOV

CREST (ENSAE), 3 avenue Pierre Larousse, 92 245 Malakoff Cedex, France;

eric.gautier@ensae.fr; alexandre.tsybakov@ensae.fr.

**ABSTRACT.** We propose a new method of estimation in high-dimensional linear regression model. It allows for very weak distributional assumptions including heteroscedasticity, and does not require the knowledge of the variance of random errors. The method is based on linear programming only, so that its numerical implementation is faster than for previously known techniques using conic programs, and it allows one to deal with higher dimensional models. We provide upper bounds for estimation and prediction errors of the proposed estimator showing that it achieves the same rate as in the more restrictive situation of fixed design and i.i.d. Gaussian errors with known variance. Following Gautier and Tsybakov (2011), we obtain the results under weaker sensitivity assumptions than the restricted eigenvalue or assimilated conditions.

## 1. INTRODUCTION

In this paper, we consider the linear regression model

$$(1) \quad y_i = x_i^T \beta^* + u_i, \quad i = 1, \dots, n,$$

where  $x_i$  are random vectors of explanatory variables in  $\mathbb{R}^p$ , and  $u_i \in \mathbb{R}$  is a random error. The aim is to estimate the vector  $\beta^* \in \mathbb{R}^p$  from  $n$  independent, not necessarily identically distributed realizations  $(y_i, x_i^T)$ ,  $i = 1, \dots, n$ . We are mainly interested in high-dimensional models where  $p$  can be much larger than  $n$  under the sparsity scenario where only few components  $\beta_k^*$  of  $\beta^*$  are non-zero ( $\beta^*$  is sparse).

The most studied techniques for high-dimensional regression under the sparsity scenario are the Lasso, the Dantzig selector, see, e.g., Candès and Tao (2007), Bickel, Ritov and Tsybakov (2009) (more references can be found in Bühlmann and van de Geer (2011) and Koltchinskii (2011)), and agregation by exponential weighting (see Dalalyan and Tsybakov (2008), Rigollet and Tsybakov (2011, 2012) and the references cited therein). Most of the literature on high-dimensional regression assumes that the random errors are Gaussian or subgaussian with known variance (or noise level). However,

quite recently several methods have been proposed which are independent of the noise level (see, e.g., Städler, Bühlmann and van de Geer (2010), Antoniadis (2010), Belloni, Chernozhukov and Wang (2011a, 2011b), Gautier and Tsybakov (2011), Sun and Zhang (2011), Belloni, Chen, Chernozhukov, and Hansen (2012) and Dalalyan (2012)). Among these, the methods of Belloni, Chernozhukov and Wang (2011b), Belloni, Chen, Chernozhukov, and Hansen (2012), Gautier and Tsybakov (2011) allow to handle non-identically distributed errors  $u_i$  and are *pivotal*, i.e., rely on very weak distributional assumptions. In Gautier and Tsybakov (2011), the regressors  $x_i$  can be correlated with the errors  $u_i$ , and an estimator is suggested that makes use of instrumental variables, called the *STIV* (Self-Tuned Instrumental Variables) estimator. In a particular instance, the *STIV* estimator can be applied in classical linear regression model where all regressors are uncorrelated with the errors. This yields a pivotal extension of the Dantzig selector based on conic programming. Gautier and Tsybakov (2011) also present a method to obtain finite sample confidence sets that are robust to non-Gaussian and heteroscedastic errors.

Another important issue is to relax the assumptions on the model under which the validity of the Lasso type methods is proved, such as the restricted eigenvalue condition of Bickel, Ritov and Tsybakov (2009) and its various analogs. Belloni, Chernozhukov and Wang (2011b) obtain fast rates for prediction for the Square-root Lasso under a relaxed version of the restricted eigenvalue condition. In the context of known noise variance, Ye and Zhang (2011) introduce cone invertibility factors instead of restricted eigenvalues. For pivotal estimation, an approach based on the sensitivities and sparsity certificates is introduced in Gautier and Tsybakov (2011), see more details below. Finally, note that aggregation by exponential weighting (Dalalyan and Tsybakov (2008), Rigollet and Tsybakov (2011, 2012)) does not require any condition on the model but its numerical realization is based on MCMC algorithms in high dimension whose convergence rate is hard to assess theoretically.

In this paper, we introduce a new pivotal estimator, called the Self-tuned Dantzig estimator. It is defined as a linear program, so from the numerical point of view it is simpler than the previously known pivotal estimators based on conic programming. We obtain upper bounds on its estimation and prediction errors under weak assumptions on the model and on the distribution of the errors showing that it achieves the same rate as in the more restrictive situation of fixed design and i.i.d. Gaussian errors with known variance. The model assumptions are based on the sensitivity analysis from Gautier and Tsybakov (2011). Distributional assumptions allow for dependence between  $x_i$  and  $u_i$ . When  $x_i$ 's are independent from  $u_i$ 's, it is enough to assume, for example, that the errors  $u_i$  are symmetric and have a finite second moment.

## 2. NOTATION

We set  $\mathbf{Y} = (y_1, \dots, y_n)^T$ ,  $\mathbf{U} = (u_1, \dots, u_n)^T$ , and we denote by  $\mathbf{X}$  the matrix of dimension  $n \times p$  with rows  $x_i^T$ ,  $i = 1, \dots, n$ . We denote by  $\mathbf{D}$  the  $p \times p$  diagonal normalizing matrix with diagonal entries  $d_{kk} > 0$ ,  $k = 1, \dots, p$ . Typical examples are:  $d_{kk} \equiv 1$  or

$$d_{kk} = \left( \frac{1}{n} \sum_{i=1}^n x_{ki}^2 \right)^{-1/2}, \quad \text{and} \quad d_{kk} = \left( \max_{i=1, \dots, n} |x_{ki}| \right)^{-1}$$

where  $x_{ki}$  is the  $k$ th component of  $x_i$ . For a vector  $\beta \in \mathbb{R}^p$ , let  $J(\beta) = \{k \in \{1, \dots, p\} : \beta_k \neq 0\}$  be its support, i.e., the set of indices corresponding to its non-zero components  $\beta_k$ . We denote by  $|J|$  the cardinality of a set  $J \subseteq \{1, \dots, p\}$  and by  $J^c$  its complement:  $J^c = \{1, \dots, p\} \setminus J$ . The  $\ell_p$  norm of a vector  $\Delta$  is denoted by  $|\Delta|_p$ ,  $1 \leq p \leq \infty$ . For  $\Delta = (\Delta_1, \dots, \Delta_p)^T \in \mathbb{R}^p$  and a set of indices  $J \subseteq \{1, \dots, p\}$ , we consider  $\Delta_J \triangleq (\Delta_1 \mathbb{1}_{\{1 \in J\}}, \dots, \Delta_p \mathbb{1}_{\{p \in J\}})^T$ , where  $\mathbb{1}_{\{\cdot\}}$  is the indicator function. For  $a \in \mathbb{R}$ , we set  $a_+ \triangleq \max(0, a)$ ,  $a_+^{-1} \triangleq (a_+)^{-1}$ .

## 3. THE ESTIMATOR

We say that a pair  $(\beta, \sigma) \in \mathbb{R}^p \times \mathbb{R}^+$  satisfies the *Self-tuned Dantzig-constraint* if it belongs to the set

$$(2) \quad \widehat{\mathcal{D}} \triangleq \left\{ (\beta, \sigma) \mid \beta \in \mathbb{R}^p, \sigma > 0, \left| \frac{1}{n} \mathbf{D} \mathbf{X}^T (\mathbf{Y} - \mathbf{X} \beta) \right|_\infty \leq \sigma r \right\}$$

for some  $r > 0$  (specified below).

**Definition 3.1.** We call the Self-Tuned Dantzig estimator any solution  $(\widehat{\beta}, \widehat{\sigma})$  of the following minimization problem

$$(3) \quad \min_{(\beta, \sigma) \in \widehat{\mathcal{D}}} \left( |\mathbf{D}^{-1} \beta|_1 + c \sigma \right),$$

for some positive constant  $c$ .

Finding the Self-Tuned Dantzig estimator is a linear program. The term  $c \sigma$  is included in the criterion to prevent from choosing  $\sigma$  arbitrarily large. The choice of the constant  $c$  will be discussed later.

## 4. SENSITIVITY CHARACTERISTICS

The sensitivity characteristics are defined by the action of the matrix

$$\Psi_n \triangleq \frac{1}{n} \mathbf{D} \mathbf{X}^T \mathbf{X} \mathbf{D}$$

on the so-called *cone of dominant coordinates*

$$C_J^{(\gamma)} \triangleq \{\Delta \in \mathbb{R}^p : |\Delta_{J^c}|_1 \leq (1 + \gamma) |\Delta_J|_1\},$$

for some  $\gamma > 0$ . It is straightforward that for  $\delta \in C_J^{(\gamma)}$ ,

$$(4) \quad |\Delta|_1 \leq (2 + \gamma) |\Delta_J|_1 \leq (2 + \gamma) |J|^{1-1/q} |\Delta_J|_q, \quad \forall 1 \leq q \leq \infty.$$

We now recall some definitions from Gautier and Tsybakov (2011). For  $q \in [1, \infty]$ , we define the  $\ell_q$  *sensitivity* as the following random variable

$$\kappa_{q,J}^{(\gamma)} \triangleq \inf_{\Delta \in C_J^{(\gamma)} : |\Delta|_q=1} |\Psi_n \Delta|_\infty.$$

Given a subset  $J_0 \subset \{1, \dots, p\}$  and  $q \in [1, \infty]$ , we define the  $\ell_q$ - $J_0$ -*block sensitivity* as

$$(5) \quad \kappa_{q,J_0,J}^{(\gamma)} \triangleq \inf_{\Delta \in C_J^{(\gamma)} : |\Delta_{J_0}|_q=1} |\Psi_n \Delta|_\infty.$$

By convention, we set  $\kappa_{q,\emptyset,J}^{(\gamma)} = \infty$ . Also, recall that the restricted eigenvalue of Bickel, Ritov and Tsybakov (2009) is defined by

$$\kappa_{\text{RE},J}^{(\gamma)} \triangleq \inf_{\Delta \in \mathbb{R}^p \setminus \{0\} : \Delta \in C_J^{(\gamma)}} \frac{|\Delta^T \Psi_n \Delta|}{|\Delta_J|_2^2}$$

and a closely related quantity is

$$\kappa'_{\text{RE},J}^{(\gamma)} \triangleq \inf_{\Delta \in \mathbb{R}^p \setminus \{0\} : \Delta \in C_J^{(\gamma)}} \frac{|J| |\Delta^T \Psi_n \Delta|}{|\Delta_J|_1^2}.$$

The next result establishes a relation between restricted eigenvalues and sensitivities. It follows directly from the Cauchy-Schwarz inequality and (4).

**Lemma 4.1.**

$$(6) \quad \kappa_{\text{RE},J}^{(\gamma)} \leq \kappa'_{\text{RE},J}^{(\gamma)} \leq (2 + \gamma) |J| \kappa_{1,J}^{(\gamma)} \leq (2 + \gamma)^2 |J| \kappa_{1,J}^{(\gamma)}.$$

The following proposition gives a useful lower bound on the sensitivity.

**Proposition 4.2.** *If  $|J| \leq s$ ,*

$$(7) \quad \kappa_{1,J,J}^{(\gamma)} \geq \frac{1}{s} \min_{k=1,\dots,p} \left\{ \min_{\Delta_k=1, |\Delta|_1 \leq (2+\gamma)s} |\Psi_n \Delta|_\infty \right\} \triangleq \kappa_{1,0}^{(\gamma)}(s).$$

**Proof.** We have

$$\begin{aligned} \kappa_{1,J,J}^{(\gamma)} &= \inf_{\Delta: |\Delta_J|_1=1, |\Delta_{J^c}|_1 \leq 1+\gamma} |\Psi_n \Delta|_\infty \\ &\geq \inf_{\Delta: |\Delta|_\infty \geq \frac{1}{s}, |\Delta|_1 \leq 2+\gamma} |\Psi_n \Delta|_\infty \\ &= \frac{1}{s} \inf_{\Delta: |\Delta|_\infty \geq 1, |\Delta|_1 \leq (2+\gamma)s} |\Psi_n \Delta|_\infty \quad (\text{by homogeneity}) \\ &= \frac{1}{s} \inf_{\Delta: |\Delta|_\infty \geq 1, |\Delta|_1 \leq (2+\gamma)s} |\Delta|_\infty \frac{|\Psi_n \Delta|_\infty}{|\Delta|_\infty} \\ &\geq \frac{1}{s} \inf_{\Delta: |\Delta|_\infty=1, |\Delta|_1 \leq (2+\gamma)s} |\Psi_n \Delta|_\infty \quad (\text{by homogeneity}) \\ &= \frac{1}{s} \inf_{\Delta: |\Delta|_\infty=1, |\Delta|_1 \leq (2+\gamma)s} |\Psi_n \Delta|_\infty \\ &= \frac{1}{s} \min_{k=1,\dots,p} \left\{ \inf_{\Delta: \Delta_k=1, |\Delta|_1 \leq (2+\gamma)s} |\Psi_n \Delta|_\infty \right\}. \quad \square \end{aligned}$$

Note that the random variable  $\kappa_{1,0}^{(\gamma)}(s)$  depends only on the observed data. It is not difficult to see that it can be obtained by solving  $p$  linear programs. For more details and further results on the sensitivity characteristics, see Gautier and Tsybakov (2011).

## 5. BOUNDS ON THE ESTIMATION AND PREDICTION ERRORS

In this section, we use the notation  $\Delta \triangleq \mathbf{D}^{-1}(\widehat{\beta} - \beta)$ . Let  $0 < \alpha < 1$  be a given constant. We choose the tuning parameter  $r$  in the definition of  $\widehat{\mathcal{D}}$  as follows:

$$(8) \quad r = \sqrt{\frac{2 \log(4p/\alpha)}{n}}.$$

**Theorem 5.1.** *Let for all  $i = 1, \dots, n$ , and  $k = 1, \dots, p$ , the random variables  $x_{ki}u_i$  be symmetric. Let  $Q^* > 0$  be a constant such that*

$$(9) \quad \mathbb{P} \left( \max_{k=1,\dots,p} \frac{d_{kk}^2}{n} \sum_{i=1}^n x_{ki}^2 u_i^2 > Q^* \right) \leq \alpha/2.$$

*Assume that  $|J(\beta^*)| \leq s$ , and set in (3)*

$$(10) \quad c = \frac{(2\gamma + 1)r}{\kappa_{1,0}^{(\gamma)}(s)},$$

where  $\gamma$  is a positive number. Then, with probability at least  $1 - \alpha$ , for any  $\gamma > 0$  and any  $\hat{\beta}$  such that  $(\hat{\beta}, \hat{\sigma})$  is a solution of the minimization problem (3) with  $c$  defined in (10) we have the following bounds on the  $\ell_1$  estimation error and on the prediction error:

$$(11) \quad |\Delta|_1 \leq \left( \frac{(\gamma + 2)(2\gamma + 1)\sqrt{Q^*}}{\gamma \kappa_{1,0}^{(\gamma)}(s)} \right) r,$$

$$(12) \quad \Delta^T \Psi_n \Delta \leq \left( \frac{(\gamma + 2)(2\gamma + 1)^2 Q^*}{\gamma^2 \kappa_{1,0}^{(\gamma)}(s)} \right) r^2.$$

**Proof.** Set

$$\hat{Q}(\beta) \triangleq \max_{k=1, \dots, p} \frac{d_{kk}^2}{n} \sum_{i=1}^n x_{ki}^2 (y_i - x_i^T \beta)^2,$$

and define the event

$$\mathcal{G} = \left\{ \left| \frac{1}{n} \mathbf{D} \mathbf{X}^T \mathbf{U} \right|_{\infty} \leq r \sqrt{\hat{Q}(\beta^*)} \right\} = \left\{ \left| \frac{d_{kk}}{n} \sum_{i=1}^n x_{ki} u_i \right| \leq r \sqrt{\hat{Q}(\beta^*)}, \quad k = 1, \dots, p \right\}.$$

Then

$$\mathcal{G}^c \subset \bigcup_{k=1, \dots, p} \left\{ \left| \frac{\sum_{i=1}^n x_{ki} u_i}{\sqrt{\sum_{i=1}^n (x_{ki} u_i)^2}} \right| \geq \sqrt{nr} \right\}$$

and the union bound yields

$$(13) \quad \mathbb{P}(\mathcal{G}^c) \leq \sum_{k=1}^p \mathbb{P} \left( \left| \frac{\sum_{i=1}^n x_{ki} u_i}{\sqrt{\sum_{i=1}^n (x_{ki} u_i)^2}} \right| \geq \sqrt{nr} \right).$$

We now use the following result on deviations of self-normalized sums due to Efron (1969).

**Lemma 5.2.** *If  $\eta_1, \dots, \eta_n$  are independent symmetric random variables, then*

$$\mathbb{P} \left( \frac{\left| \frac{1}{n} \sum_{i=1}^n \eta_i \right|}{\sqrt{\frac{1}{n} \sum_{i=1}^n \eta_i^2}} \geq t \right) \leq 2 \exp \left( -\frac{nt^2}{2} \right), \quad \forall t > 0.$$

For each of the probabilities on the right-hand side of (13), we apply Lemma 5.2 with  $\eta_i = x_{ki} u_i$ . This and the definition of  $r$  yield  $\mathbb{P}(\mathcal{G}^c) \leq \alpha/2$ . Thus, the event  $\mathcal{G}$  holds with probability at least  $1 - \alpha/2$ . On the event  $\mathcal{G}$  we have

$$(14) \quad |\Psi_n \Delta|_{\infty} \leq \left| \frac{1}{n} \mathbf{D} \mathbf{X}^T (\mathbf{Y} - \mathbf{X} \hat{\beta}) \right|_{\infty} + \left| \frac{1}{n} \mathbf{D} \mathbf{X}^T (\mathbf{Y} - \mathbf{X} \beta^*) \right|_{\infty}$$

$$(15) \quad \begin{aligned} &\leq r \hat{\sigma} + \left| \frac{1}{n} \mathbf{D} \mathbf{X}^T \mathbf{U} \right|_{\infty} \\ &\leq r \left( \hat{\sigma} + \sqrt{\hat{Q}(\beta^*)} \right) \end{aligned}$$

$$(16) \quad \leq r \left[ 2\sqrt{\widehat{Q}(\beta^*)} + \left( \widehat{\sigma} - \sqrt{\widehat{Q}(\beta^*)} \right) \right]$$

Inequality (15) holds because  $(\widehat{\beta}, \widehat{\sigma})$  belongs to the set  $\widehat{\mathcal{D}}$  by definition. Notice that, on the event  $\mathcal{G}$ ,  $(\beta^*, \sqrt{\widehat{Q}(\beta^*)})$  belongs to the set  $\widehat{\mathcal{D}}$ . On the other hand,  $(\widehat{\beta}, \widehat{\sigma})$  minimizes the criterion  $|\mathbf{D}^{-1}\beta|_1 + c\sigma$  on the same set  $\widehat{\mathcal{D}}$ . Thus, on the event  $\mathcal{G}$ ,

$$(17) \quad |\mathbf{D}^{-1}\widehat{\beta}|_1 + c\widehat{\sigma} \leq |\mathbf{D}^{-1}\beta^*|_1 + c\sqrt{\widehat{Q}(\beta^*)}.$$

This implies, again on the event  $\mathcal{G}$ ,

$$(18) \quad \begin{aligned} |\Psi_n \Delta|_\infty &\leq r \left[ 2\sqrt{\widehat{Q}(\beta^*)} + \frac{1}{c} \sum_{k \in J(\beta^*)} \left( |d_{kk}^{-1}\beta_k^*| - |d_{kk}^{-1}\widehat{\beta}_k| \right) - \frac{1}{c} \sum_{k \in J(\beta^*)^c} |d_{kk}^{-1}\widehat{\beta}_k| \right] \\ &\leq r \left( 2\sqrt{\widehat{Q}(\beta^*)} + \frac{1}{c} |\Delta_{J(\beta^*)}|_1 \right) \end{aligned}$$

where  $\beta_k^*, \widehat{\beta}_k$  are the  $k$ th components of  $\beta^*, \widehat{\beta}$ . Similarly, (17) implies that, on the event  $\mathcal{G}$ ,

$$(19) \quad \begin{aligned} |\Delta_{J(\beta^*)^c}|_1 &= \sum_{k \in J(\beta^*)^c} |d_{kk}^{-1}\widehat{\beta}_k| \\ &\leq \sum_{k \in J(\beta^*)} \left( |d_{kk}^{-1}\beta_k^*| - |d_{kk}^{-1}\widehat{\beta}_k| \right) + c \left( \sqrt{\widehat{Q}(\beta^*)} - \widehat{\sigma} \right) \\ &\leq |\Delta_{J(\beta^*)}|_1 + c\sqrt{\widehat{Q}(\beta^*)}. \end{aligned}$$

We now distinguish between the following two cases.

Case 1:  $c\sqrt{\widehat{Q}(\beta^*)} \leq \gamma |\Delta_{J(\beta^*)}|_1$ . In this case (19) implies

$$(20) \quad |\Delta_{J(\beta^*)^c}|_1 \leq (1 + \gamma) |\Delta_{J(\beta^*)}|_1.$$

Thus,  $\Delta \in C_{J(\beta^*)}^{(\gamma)}$  on the event  $\mathcal{G}$ . By definition of  $\kappa_{1, J(\beta^*), J(\beta^*)}^{(\gamma)}$  and (7),

$$|\Delta_{J(\beta^*)}|_1 \leq \frac{|\Psi_n \Delta|_\infty}{\kappa_{1, J(\beta^*), J(\beta^*)}^{(\gamma)}} \leq \frac{|\Psi_n \Delta|_\infty}{\kappa_{1, 0}^{(\gamma)}(s)}.$$

This and (18) yield

$$|\Delta_{J(\beta^*)}|_1 \leq \frac{2r\sqrt{\widehat{Q}(\beta^*)}}{\kappa_{1, 0}^{(\gamma)}(s)} \left( 1 - \frac{r}{c\kappa_{1, 0}^{(\gamma)}(s)} \right)_+^{-1}.$$

Case 2:  $c\sqrt{\widehat{Q}(\beta^*)} > \gamma |\Delta_{J(\beta^*)}|_1$ . Then, obviously,  $|\Delta_{J(\beta^*)}|_1 < \frac{c}{\gamma} \sqrt{\widehat{Q}(\beta^*)}$ .

Combining the two cases we obtain, on the event  $\mathcal{G}$ ,

$$(21) \quad |\Delta_{J(\beta^*)}|_1 \leq \sqrt{\widehat{Q}(\beta^*)} \max \left\{ \frac{2r}{\kappa_{1,0}^{(\gamma)}(s)} \left( 1 - \frac{r}{c\kappa_{1,0}^{(\gamma)}(s)} \right)_+^{-1}, \frac{c}{\gamma} \right\}.$$

In this argument,  $c > 0$  and  $\gamma > 0$  were arbitrary. The value of  $c$  given in (10) is the minimizer of the right-hand side of (21). Plugging it in (21) we find that, with probability at least  $1 - \alpha/2$

$$|\Delta|_1 \leq \frac{(\gamma + 2)(2\gamma + 1)r}{\gamma\kappa_{1,0}^{(\gamma)}(s)} \sqrt{\widehat{Q}(\beta^*)}$$

where we have used (19). Now, by (9),  $\widehat{Q}(\beta^*) \leq Q^*$  with probability at least  $1 - \alpha/2$ . Thus, we get that (11) holds with probability at least  $1 - \alpha$ . Next, using (18) we obtain that, on the same event of probability at least  $1 - \alpha$ ,

$$|\Psi_n \Delta|_\infty \leq \frac{(2\gamma + 1)r}{\gamma} \sqrt{Q^*}.$$

Combining this inequality with (11) yields (12).  $\square$

#### *Discussion of Theorem 5.1.*

- (1) In view of Lemma 4.1,  $\kappa_{1,J(\beta^*),J(\beta^*)}^{(\gamma)} \geq (2 + \gamma)^{-2} \kappa_{\text{RE},J(\beta^*)}^{(\gamma)} / s$ . Also, it is easy to see from Proposition 4.2 that  $\kappa_{1,0}^{(\gamma)}(s)$  is of the order  $1/s$  when  $\Psi_n$  is the identity matrix and  $p \gg s$  (this is preserved for  $\Psi_n$  that are small perturbations of the identity). Thus, the bounds (11) and (12) take the form

$$|\Delta|_1 \leq C \left( s \sqrt{\frac{\log p}{n}} \right), \quad \Delta^T \Psi_n \Delta \leq C \left( \frac{s \log p}{n} \right),$$

for some constant  $C$ , and we recover the usual rates for the  $\ell_1$  estimation and for the prediction error respectively, cf. Bickel, Ritov and Tsybakov (2009).

- (2) Theorem 5.1 does not assume that  $x_{ki}$ 's are independent from  $u_i$ 's. The only assumption is the symmetry of  $x_{ki}u_i$ . However, if  $x_{ki}$  is independent from  $u_i$ , then by conditioning on  $x_{ki}$  in the bound for  $\mathbb{P}(\mathcal{G})$ , it is enough to assume the symmetry of  $u_i$ 's. Furthermore, while we have chosen the symmetry since it makes the conditions of Theorem 5.1 simple and transparent, it is not essential for our argument to be applied. The only point in the proof where we use the symmetry is the bound for the probability of deviations of self-normaized sums  $\mathbb{P}(\mathcal{G})$ . This probability can be bounded in many other ways without the symmetry assumption, cf., e.g., Gautier and Tsybakov (2011). It is enough to have  $\mathbb{E}[x_{ki}u_i] = 0$  and a uniform over  $k$  control



of the ratio

$$\frac{(\sum_{i=1}^n \mathbb{E}[x_{ki}^2 u_i^2])^{1/2}}{(\sum_{i=1}^n \mathbb{E}[|x_{ki} u_i|^{2+\delta}])^{1/(2+\delta)}}$$

for some  $\delta > 0$ , cf. [14] or [6].

- (3) The quantity  $Q^*$  is not present in the definition of the estimator and is needed only to assess the rate of convergence. It is not hard to find  $Q^*$  in various situations. The simplest case is when  $d_{kk} \equiv 1$  and the random variables  $x_{ki}$  and  $u_i$  are bounded uniformly in  $k, i$  by a constant  $L$ . Then we can take  $Q^* = L^4$ . If only  $x_{ki}$  are bounded uniformly in  $k$  by  $L$ , condition (9) holds when  $\mathbb{P}(\frac{1}{n} \sum_{i=1}^n u_i^2 > Q^*/L^2) \leq \alpha/2$ , and then for  $Q^*$  to be bounded it is enough to assume that  $u_i$ 's have a finite second moment. The same remark applies when  $d_{kk} = (\max_{i=1, \dots, n} |x_{ki}|)^{-1}$ , with an advantage that in this case we guarantee that  $Q^*$  is bounded under no assumption on  $x_{ki}$ .
- (4) The bounds in Theorem 5.1 depend on  $\gamma > 0$  that can be optimized. Indeed, the functions of  $\gamma$  on the right-hand sides of (11) and (12) are data-driven and can be minimized on a grid of values of  $\gamma$ . Thus, we obtain an optimal (random) value  $\gamma = \hat{\gamma}$ , for which (11) and (12) remain valid, since these results hold for any  $\gamma > 0$ .

## REFERENCES

- [1] Antoniadis, A. (2010) Comments on:  $l^1$ -penalization for Mixture Regression Models. (with discussion). *Test*, **19**, 257–258.
- [2] Belloni, A., Chen, D., Chernozhukov, V. and Hansen C. (2012) Sparse Models and Methods for Optimal Instruments with an Application to Eminent Domain. *Econometrica*, **80**, 2369–2430.
- [3] Belloni, A. and Chernozhukov, V. (2010) Least Squares After Model Selection in High-dimensional Sparse Models. Forthcoming in *Bernoulli*.
- [4] Belloni, A., Chernozhukov, V. and Wang, L. (2011a) Square-Root Lasso: Pivotal Recovery of Sparse Signals via Conic Programming. *Biometrika*, **98**, 791–806.
- [5] Belloni, A., Chernozhukov, V. and Wang, L. (2011b) Pivotal Estimation of Nonparametric Functions via Square-root Lasso. Preprint <http://arxiv.org/pdf/1105.1475.pdf>
- [6] Bertail, P. , Gauth  rat, E. and Harari-Kermadec, H. (2009) Exponential Inequalities for Self Normalized Sums. *Electronic Communications in Probability*, **13**, 628–640.
- [7] Bickel, P., Ritov, J. Y. and Tsybakov, A. B. (2009) Simultaneous Analysis of Lasso and Dantzig Selector. *The Annals of Statistics*, **37**, 1705–1732.
- [8] B  hlmann, P. and van de Geer, S.A. (2011) *Statistics for High-Dimensional Data*. Springer, New-York.
- [9] Cand  s, E., and Tao, T. (2007) The Dantzig Selector: Statistical Estimation when  $p$  is Much Larger than  $n$ . *The Annals of Statistics*, **35**, 2313–2351.

- [10] Dalalyan, A. (2012) SOCP Based Variance Free Dantzig Selector with Application to Robust Estimation. *C. R. Math. Acad. Sci. Paris*, **350**, 785–788
- [11] Dalalyan, A., and Tsybakov, A.B. (2008) Aggregation by Exponential Weighting, Sharp PAC-Bayesian Bounds and Sparsity. *Journal of Machine Learning Research*, **72**, 39–61.
- [12] Efron, B. (1969) Student's t-test Under Symmetry Conditions. *Journal of American Statistical Association*, **64**, 1278–1302.
- [13] Gautier, E. and Tsybakov, A.B. (2011) High-dimensional instrumental variables regression and confidence sets. Preprint <http://arxiv.org/pdf/1105.2454v1.pdf>
- [14] Jing, B.-Y., Shao, Q. M. and Wang, Q. (2003) Self-Normalized Cramér-Type Large Deviations for Independent Random Variables. *The Annals of Probability*, **31**, 2167–2215.
- [15] Koltchinskii, V. (2011) *Oracle Inequalities for Empirical Risk Minimization and Sparse Recovery Problems*. Lecture Notes in Mathematics, vol. 2033. Springer, New-York.
- [16] Rigollet, P. and Tsybakov, A.B. (2011) Exponential Screening and Optimal Rates of Sparse Estimation. *The Annals of Statistics*, **39**, 731–771.
- [17] Rigollet, P. and Tsybakov, A.B. (2012) Sparse Estimation by Exponential Weighting. *Statistical Science*, **27**, 558–575.
- [18] Städler, N., Bühlmann, P. and van de Geer, S.A. (2010)  $l^1$ -penalization for Mixture Regression Models. *Test*, **19**, 209–256.
- [19] Sun, T. and Zhang, C.-H. (2011) Scaled Sparse Linear Regression. Preprint <http://arxiv.org/abs/1104.4595>
- [20] Ye, F. and Zhang, C.-H. (2010) Rate Minimality of the Lasso and Dantzig Selector for the  $l_q$  Loss in  $l_r$  Balls. *Journal of Machine Learning Research*, **11**, 3519–3540.